

STAT 378 Assignment 4 Final Project

Siwei Li Siqu Wang Ruhan Zhang

November 6, 2015

Problem 1. Metropolis-Hastings Algorithm

Step1: Main Algorithmmm for Beta distribution-Metropolis Hasting

```
set.seed(99)
```

```
#prepare for the traceplot in later stages:
```

```
library(coda)
```

```
## Warning: package 'coda' was built under R version 3.1.3
```

```
#Metropolis Hasting Procedure:
```

```
Beta_Metropolis <- function(alpha=6,beta=4,c,iterations){
```

```
  chain <- numeric(iterations+1)
```

```
  start_value <- runif(1)
```

```
  chain[1] <- start_value
```

```
  for(i in 1:iterations){
```

```
    current_value <- chain[i]
```

```
    #proposal function:
```

```
    newbeta <- 1
```

```
    while (newbeta == 0 | newbeta == 1) {
```

```
      newbeta <- rbeta(1, c*current_value, c*(1-current_value))
```

```
    }
```

```
    #proposal ratio:
```

```
    proposal_ratio <- dbeta(current_value, c*newbeta, c*(1-newbeta)) / dbeta(newbeta, c*current_value, c)
```

```
    #posterior ratio:
```

```
    posterior_ratio <- dbeta(newbeta, alpha, beta) / dbeta(current_value, alpha, beta)
```

```
    #acceptance ratio:
```

```
    if(runif(1) < min(1, posterior_ratio * proposal_ratio)){
```

```
      chain[i+1] <- newbeta
```

```
    }else{
```

```
      chain[i+1] <- current_value
```

```
    }
```

```
  }
```

```
  return(chain)
```

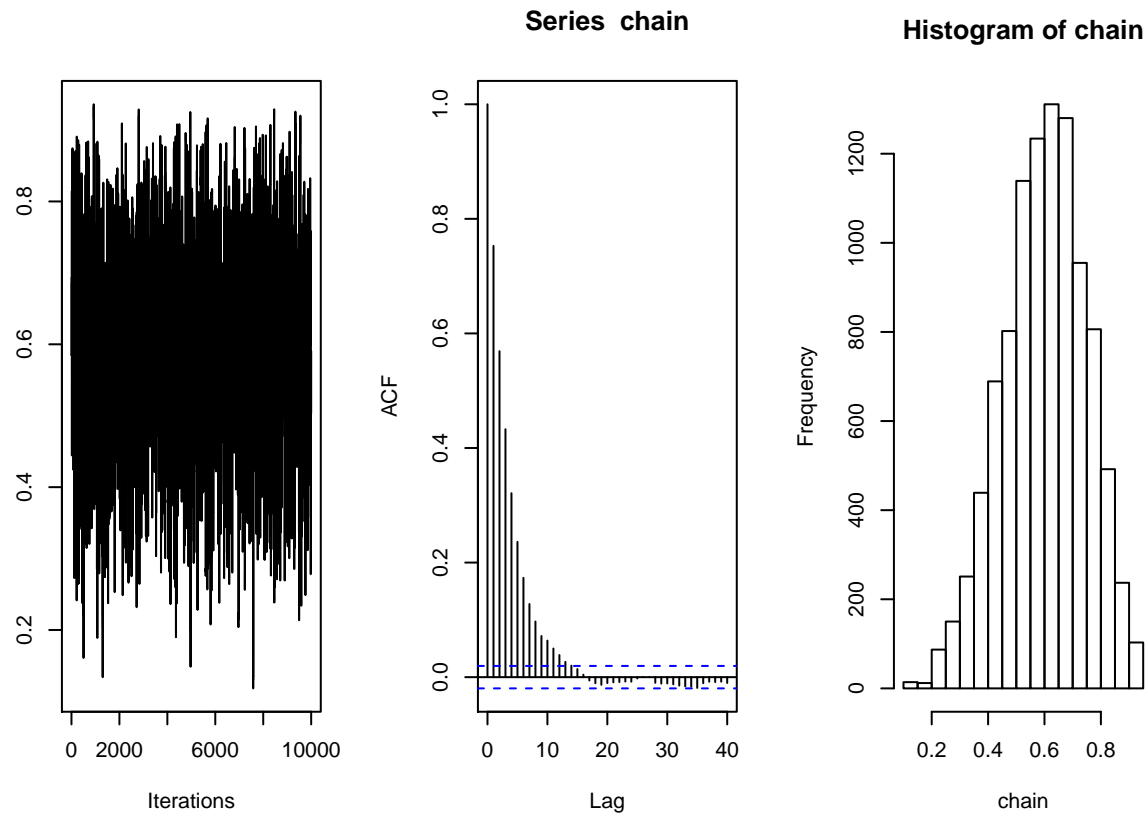
```
}
```

Step2: evaluate the performance of the sampler

```
chain <- Beta_Metropolis(alpha=6,beta=4,c=1,iterations=10000)
```

```
par(mfrow=c(1,3)) #1 row, 3 columns
```

```
traceplot(as.mcmc(chain)); acf(chain); hist(chain) #plot commands
```

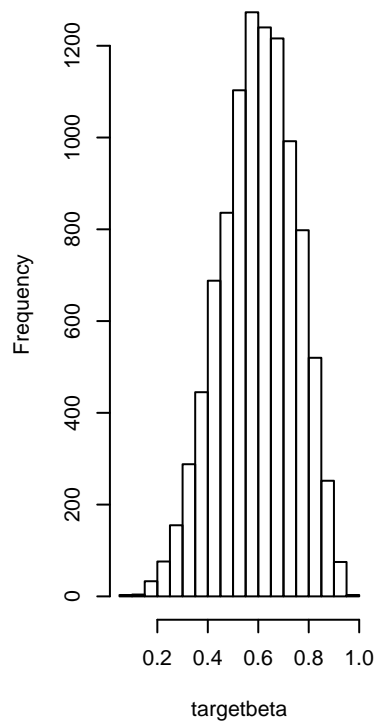


```
#graphical comparison
targetbeta <- rbeta(10000,6,4)
hist(targetbeta)
#numerical comparison - Kolmogorov-Smirnov statistic
ks.test(chain,targetbeta)
```

```
## Warning in ks.test(chain, targetbeta): p-value will be approximate in the
## presence of ties
```

```
##
## Two-sample Kolmogorov-Smirnov test
##
## data: chain and targetbeta
## D = 0.0193, p-value = 0.04914
## alternative hypothesis: two-sided
```

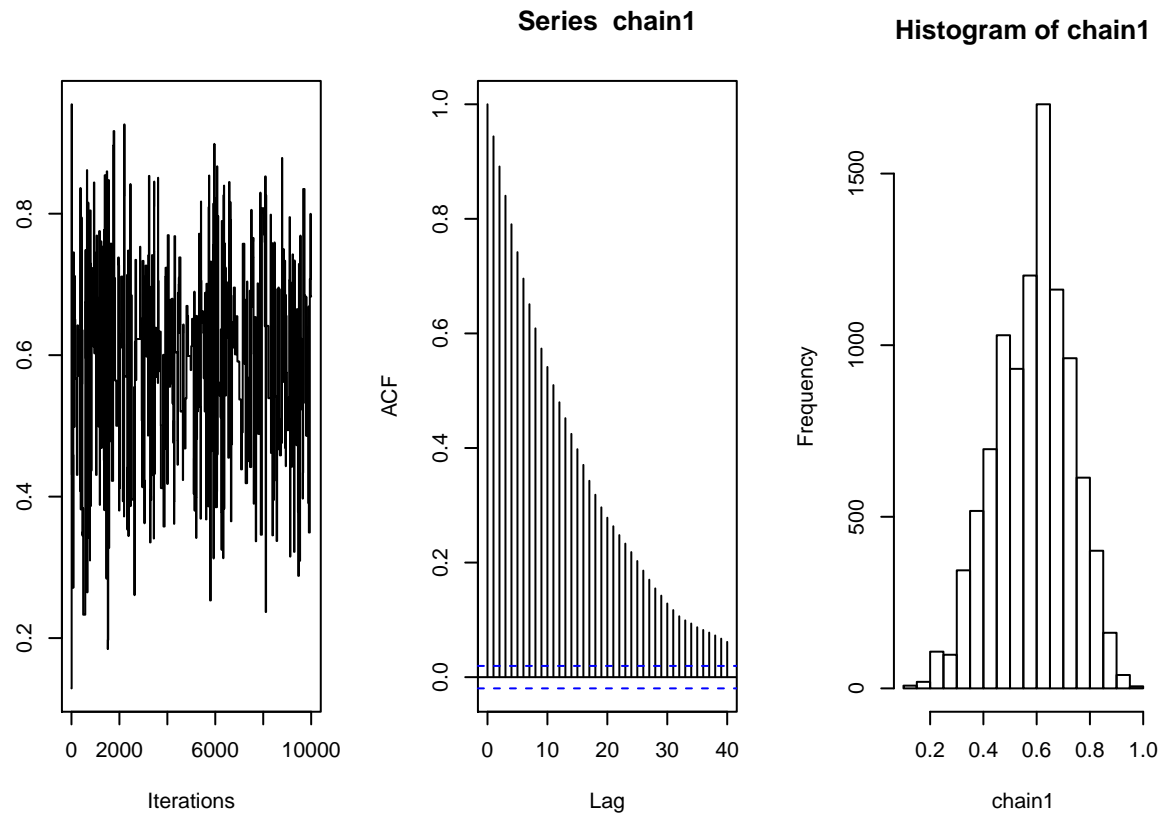
Histogram of targetbeta



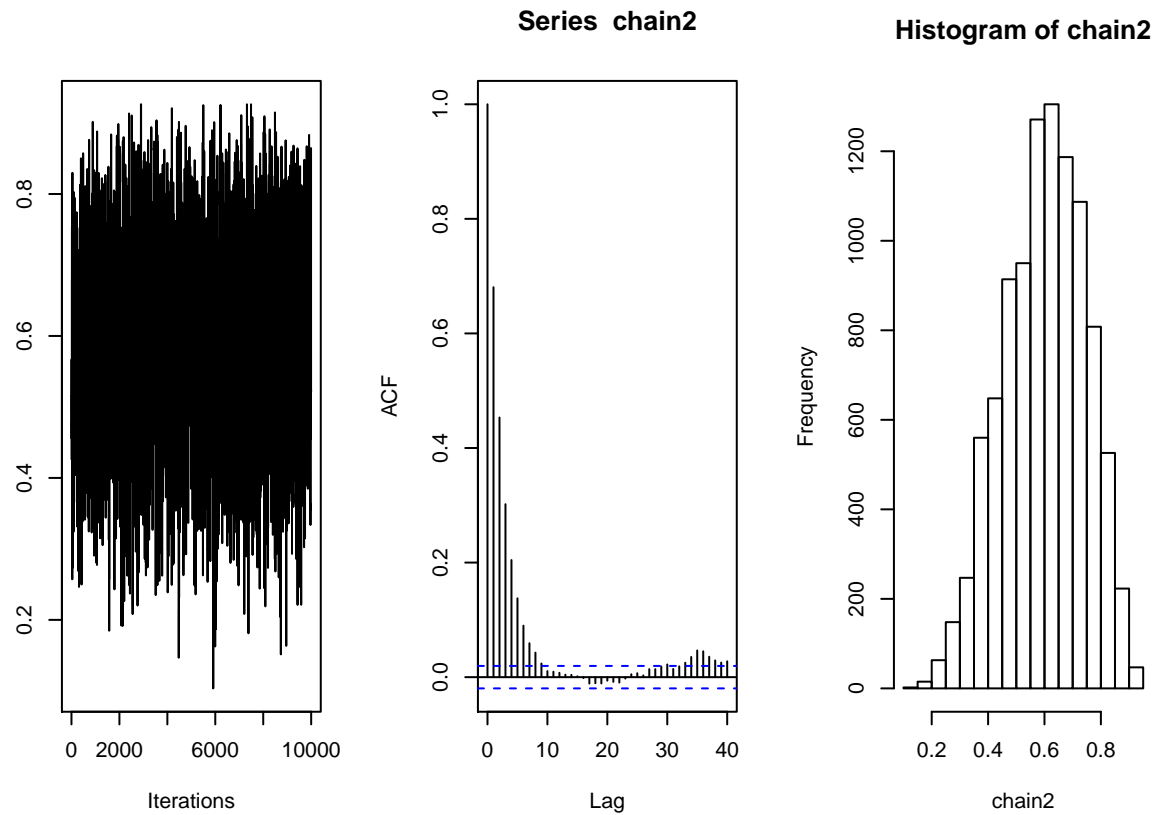
Using Kolmogorov-Smirnov test, we are trying to test the if our chain follows our target Beta distribution, $\text{Beta}(\text{'alpha'}=6, \text{'beta'}=4)$. The null hypothesis of this test states that our chain follows the target Beta distribution. The p-value for this test is $0.04914 < 0.05$ if we are evaluating at 95% significance level. Therefore, we do not reject the null hypothesis. Our metropolis hasting performs quite good in getting samples from the target Beta distribution.

Step3: re-run the sampler with $c=0.1$, $c=2.5$ and $c=10$

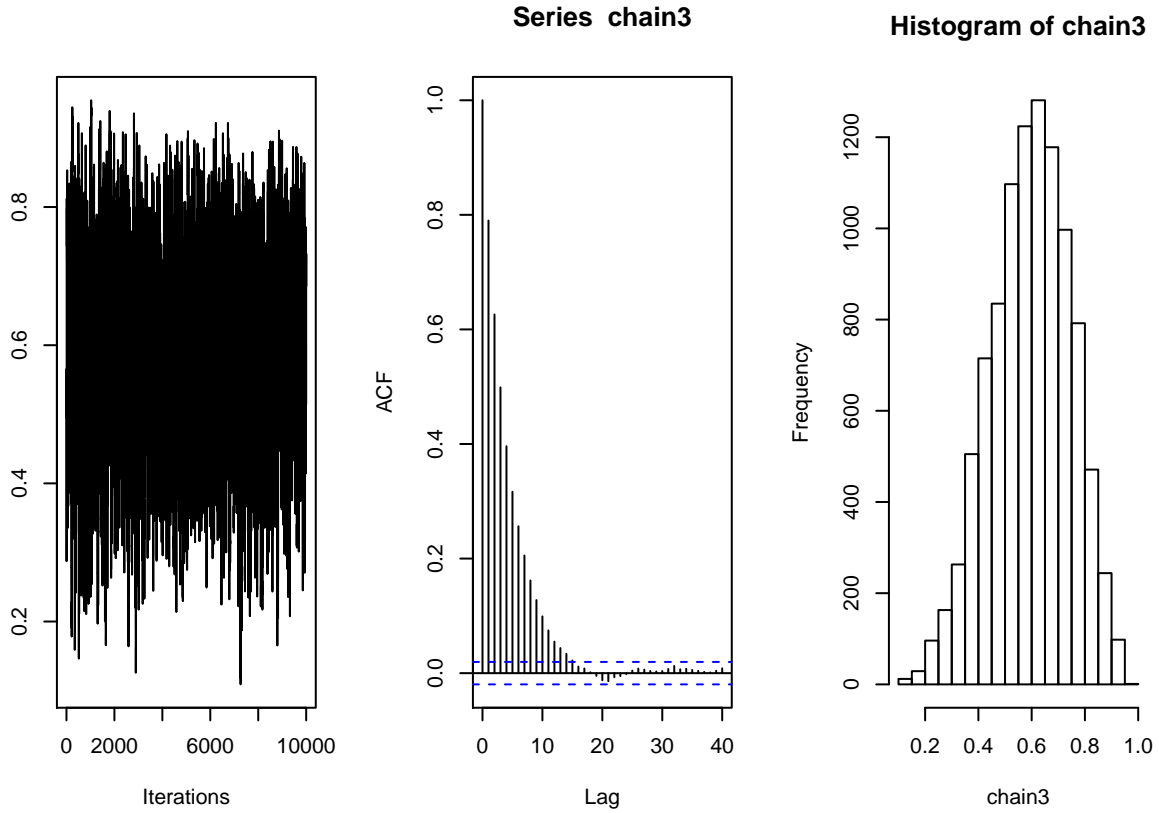
```
par(mfrow=c(1,3))
chain1 <- Beta_Metropolis(alpha=6,beta=4,c=0.1,iterations=10000)
traceplot(as.mcmc(chain1)); acf(chain1); hist(chain1)
```



```
chain2 <- Beta_Metropolis(alpha=6,beta=4,c=2.5,iterations=10000)
traceplot(as.mcmc(chain2)); acf(chain2); hist(chain2)
```



```
chain3 <- Beta_Metropolis(alpha=6,beta=4,c=10,iterations=10000)
traceplot(as.mcmc(chain3)); acf(chain3); hist(chain3)
```



Problem 2. Gibbs Sampling

To sample from the full conditional distribution, we need to figure out the full conditional density function first. The problem indicates that

$$p(x|y) \propto ye^{-yx} \quad 0 < x < B < \infty$$

$$p(y|x) \propto xe^{-yx} \quad 0 < y < B < \infty$$

Suppose the normalizing constant of $p(x|y)$ is c , we can figure out c by

$$\begin{aligned} \int_0^B cye^{-yx} dx &= 1 \\ c(-e^{-yx}|_0^B) &= 1 \\ c(1 - e^{-By}) &= 1 \end{aligned}$$

So now we know that $c = \frac{1}{1 - e^{-By}}$. Similarly, we can find that the normalizing constant for $p(y|x)$ is $\frac{1}{1 - e^{-Bx}}$. And we can write the density function as

$$p(x|y) = \frac{ye^{-yx}}{1 - e^{-By}} \quad 0 < x < B < \infty$$

$$p(y|x) = \frac{xe^{-yx}}{1 - e^{-Bx}} \quad 0 < y < B < \infty$$

Then the cumulative density functions are

$$F_{X|Y}(x|y) = \frac{1 - e^{-yx}}{1 - e^{-By}} \quad 0 < x < B$$

$$F_{Y|X}(y|x) = \frac{1 - e^{-yx}}{1 - e^{-Bx}} \quad 0 < y < B$$

With some computations, we may obtain the inverse function of two CDFs:

$$F_{X|Y}^{(-1)}(u|x, y) = \frac{-\log(1 - u + ue^{-By})}{y} \quad 0 < x < B$$

$$F_{Y|X}^{(-1)}(u|x, y) = \frac{-\log(1 - u + ue^{-Bx})}{x} \quad 0 < y < B$$

where $0 \leq u \leq 1$. Now, we can list the algorithm scheme:

1. Select initial values (x_0, y_0) and set the iteration number T .
2. Set index $t = 1$.
3. Generate random draws (x_t, y_t) where
 - Generate two numbers u_1, u_2 from Unif(0,1)
 - Let $x_t = F_{X|Y}^{(-1)}(u_1|x_{t-1}, y_{t-1})$
 - Let $y_t = F_{Y|X}^{(-1)}(u_2|x_t, y_{t-1})$
4. Record sample (x_t, y_t)
5. Set $t = t + 1$, return to step 3 if $t \leq T$.
6. Return T samples.

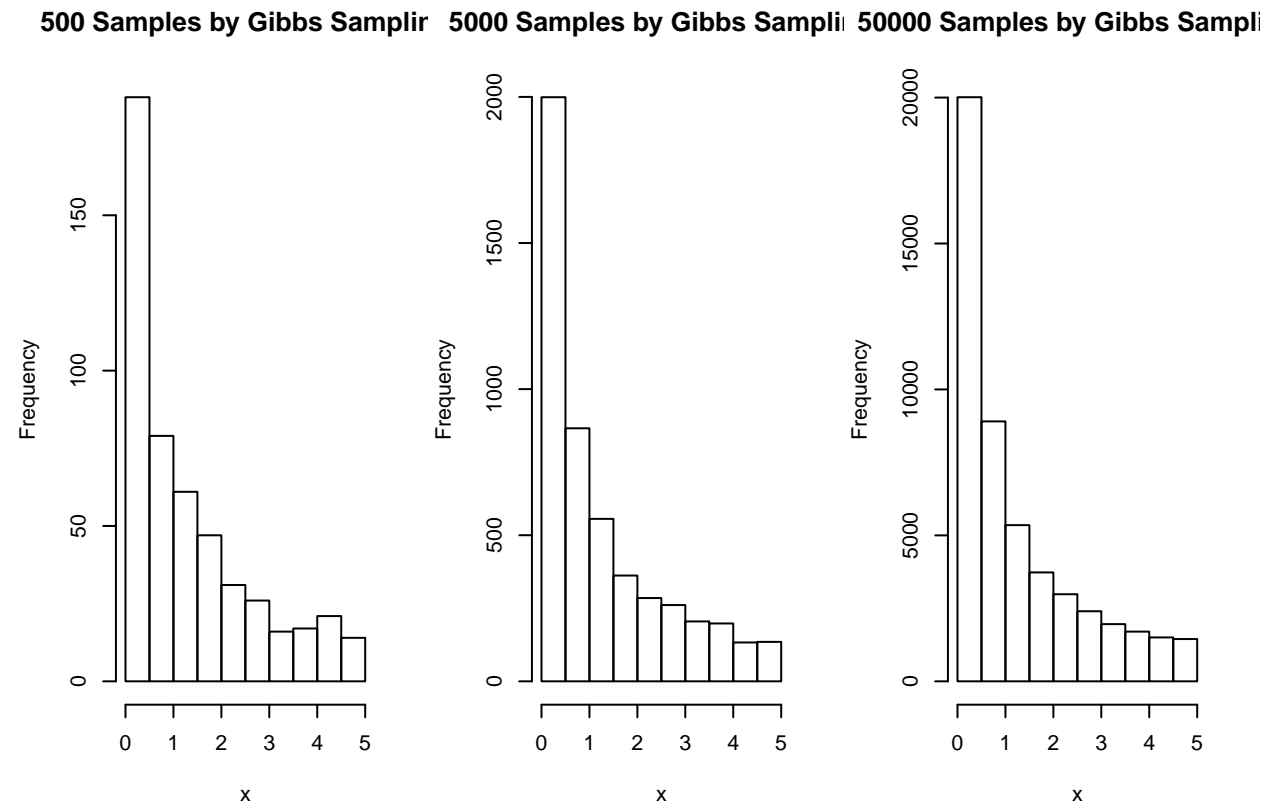
The following is the R code of the algorithm, we denote the Gibbs sampling function as Gibbs.sim():

```
Gibbs.sim <- function(x0,y0,B,Ite){#The inputs are initial values x0, y0, the bound B and the iteration
  if (x0<0 || x0>B || y0<0 || y0>B ) {#Check the validity of initial values
    stop("The initial value is not in the correct range.")
  }
  else {x_result=c()
    y_result=c()
    xt=x0
    yt=y0
    for(t in 1:Ite){
      u1=runif(1)
      u2=runif(1)
      xt=-log(1-u1*(1-exp(-B*yt)))/yt
      yt=-log(1-u2*(1-exp(-B*xt)))/xt #the xt here is the updated why, thus the algorithm construct
      x_result=c(x_result, xt)
      y_result=c(y_result, yt)
    }
    return(cbind(x_result,y_result)) #We will do burn in later
  }
}
```

We do not let the function automatically do the burn-in since we think it is not proper to discard the sampling information easily. Besides, the burn-in process is not difficult to be implemented in this case. Now for $B = 5$, we made three samples with sizes $T = 500, 5000, 50000$ and plot the histograms. For each sample, we burn-in the first 200 samples. In order to ensure the correct sample amount, we actually did 700, 5200, 50200 iterations each time.

```
set.seed(123)
sample500=Gibbs.sim(1,2,5,700)[- (1:200),] #We burn-in the first 200 samples in order to get 500 samples
sample5000=Gibbs.sim(1,2,5,5200)[- (1:200),] #Also burn-in the first 200 samples
sample50000=Gibbs.sim(1,2,5,50200)[- (1:200),] #50000 samples after burn-in

### Plot the histogram ###
par(mfrow=c(1,3))
hist(sample500[,1],xlab="x", main="500 Samples by Gibbs Sampling")
hist(sample5000[,1],xlab="x", main="5000 Samples by Gibbs Sampling")
hist(sample50000[,1],xlab="x", main="50000 Samples by Gibbs Sampling")
```



We can see that all the histograms are skewed to the right. The shape of each plot is similar. \

By Gibbs sampling, once the Markov Chain constructed by the algorithm converges to the target distribution which we want to find, we can treat samples as approximates from the target distribution. Thus we can use the mean of samples of x to estimate $E_{p(X)}[X]$. Here are the estimates from three sizes samples:

```
mean(sample500[,1])
```

```
## [1] 1.334719
```



```
mean(sample5000[,1])
```

```
## [1] 1.265094
```

```
mean(sample50000[,1])
```

```
## [1] 1.258227
```

The mean of the sample with $T = 500$ is about 1.07 and the other two means are approximately 1.25 and 1.27. We can see that mean estimates of samples with $T = 5000$ and $T = 50000$ are much closer. It is a reasonable phenomena since as the iteration number increases, we will get a more accurate approximation from the target distribution. It indicates that we may get a good estimate of $E_{p(X)}[X]$ from samples with $T > 5000$.

Problem 3. K-Means

Before starting this part, we first load two packages: fpc and flexclust. The package fpc will allow us to use the command plotcluster() to visualize the clustering result. The package flexclust will allow us to use the function randIndex() to evaluate the clustering result.

```
library('fpc')
```

```
## Warning: package 'fpc' was built under R version 3.1.3
```

```
library('flexclust')
```

```
## Warning: package 'flexclust' was built under R version 3.1.3
```

```
## Loading required package: grid  
## Loading required package: lattice  
## Loading required package: modeltools
```

```
## Warning: package 'modeltools' was built under R version 3.1.3
```

```
## Loading required package: stats4
```

Instead of using the built-in function kmeans() directly, we try to write the K means algorithm by ourselves. The function is called clusterByKmeans() with two parameter inputs: the data set and the number of groups.

```
clusterByKmeans <- function(dataset, groupNo){  
  
  # Obtain the number of rows and the number of columns of the dataset  
  rowNo <- nrow(dataset)  
  colNo <- ncol(dataset)  
  
  # Create a vector g to store the clustering result  
  # (indicating which point belongs to which group)  
  cluster <- vector("numeric", length=rowNo)
```

```

# Create a vector to store the distance between each data point
# and the center of the cluster to which that data point belongs.
dtc <- vector("numeric", length=rowNo)

# Create a matrix to store the centers of the groups
center <- matrix(0,nrow=groupNo,ncol=colNo)

# Randomly select 3 points (rows of data.train) as the initial centers of the 3 groups
randomInt <- as.vector(sample(1:rowNo,size=groupNo))
for (i in 1:groupNo){
  center[i,] <- dataset[randomInt[i],]
  center <- matrix(center,groupNo,colNo)
}

# We specify an initial way of clustering:
# all the rest of the points belong to the first cluster,
# while the points chosen as the initial centers belong to its own cluster.
for (i in 1:rowNo){
  cluster[i] <- 1
}
for (j in 1:groupNo){
  cluster[randomInt[j]] <- j
}

# Calculate the initial value of the distance vector, dtc
for (i in 1:rowNo){
  dtc[i] <- sqrt(sum((dataset[i,]-center[cluster[i],])^2))
}

# Create a variable, changed, which indicates
# whether the clustering result is changed in an iteration.
changed <- TRUE

# If the assignment is changed in a iteration, we continue;
# if the assignment no longer changes, we stop.
while(changed){

initialCluster <- cluster

for (i in 1:rowNo){
  initialdtc <- dtc[i]
  preCluster <- cluster[i]

  # We shall find a better cluster solution for the point i.
  for (j in 1:groupNo){
    # Compute the distance between the point i and the center of the cluster j
    currentdtc <- sqrt(sum((dataset[i,]-center[j,])^2))
    if (currentdtc<initialdtc){
      # Update the clustering result
      cluster[i] <- j
      # Update the distance vector
      dtc[i] <- currentdtc
    }
  }
}
}

```

```

    }
  }
}

# Compare the new assignment with the previous assignment.
# The overall assignment is considered as "changed" if the assignment of
# at least one point is changed.
for (i in 1:rowNo){
  if (cluster[i] != initialCluster[i]){
    changed <- TRUE
    break
  }else{
    changed <- FALSE
  }
}

# Then we update the centers of each cluster under the new assignment.
for (k in 1:groupNo){
  # Obtain all the points in the cluster k
  pointsIn <- dataset[cluster==k,]
  pointsMat <- as.matrix(pointsIn)
  # If there is at least one points in the cluster k, we update the center
  # of the cluster k; otherwise, the center remains unchanged.
  if (nrow(pointsMat)>0){
    center[k,] <- colMeans(pointsMat)
  }else{
    center[k,] <- center[k,]
  }
}
}
return(cluster)
}

```

Now we shall use the function `clusterByKmeans()` to cluster the wine data and the iris data.

K means for the wine data without scaling:

```

# Obtain the wine data from the package "rattle"
data(wine, package="rattle")
head(wine)

```

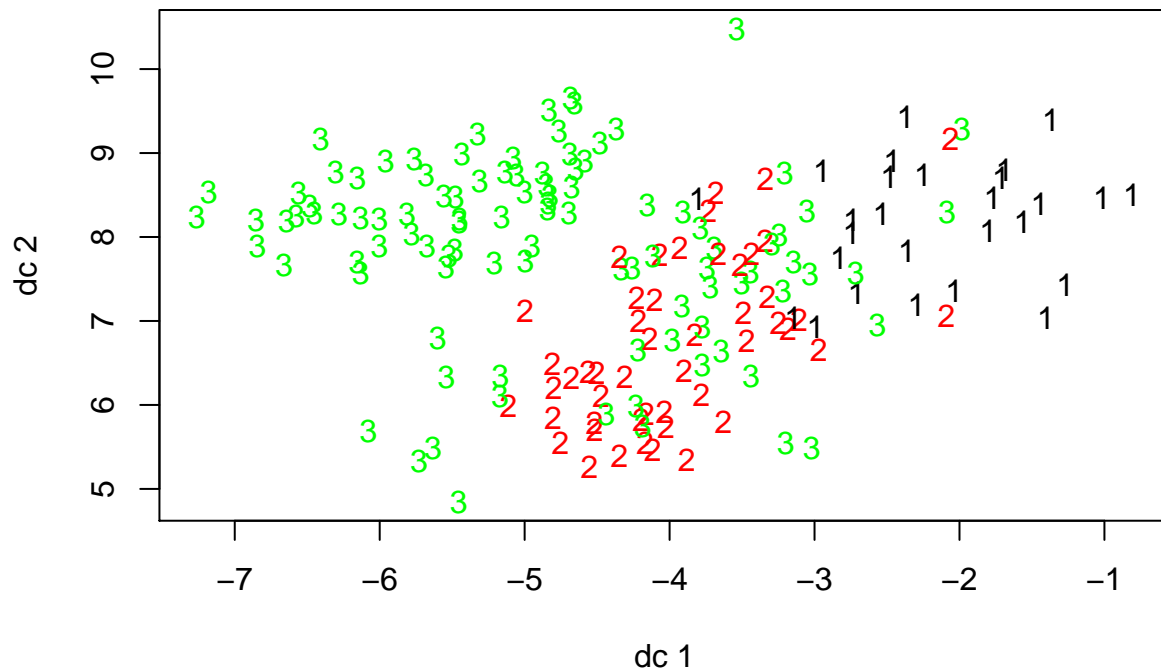
```

##   Type Alcohol Malic  Ash Alkalinity Magnesium Phenols Flavanoids
## 1    1   14.23  1.71 2.43      15.6      127    2.80      3.06
## 2    1   13.20  1.78 2.14      11.2      100    2.65      2.76
## 3    1   13.16  2.36 2.67      18.6      101    2.80      3.24
## 4    1   14.37  1.95 2.50      16.8      113    3.85      3.49
## 5    1   13.24  2.59 2.87      21.0      118    2.80      2.69
## 6    1   14.20  1.76 2.45      15.2      112    3.27      3.39
##   Nonflavanoids Proanthocyanins Color  Hue Dilution Proline
## 1           0.28           2.29 5.64 1.04      3.92    1065
## 2           0.26           1.28 4.38 1.05      3.40    1050
## 3           0.30           2.81 5.68 1.03      3.17    1185
## 4           0.24           2.18 7.80 0.86      3.45    1480
## 5           0.39           1.82 4.32 1.04      2.93     735

```

```
## 6          0.34          1.97  6.75 1.05          2.85          1450
```

```
# Exclude the "Type" variable from the data inputs
# Assign the others to data.train
data.train <- wine[-1]
# Use set seed() so that the clustering results can be reproducible
set.seed(37810)
# Use k-means method to cluster the wines into 3 groups
fit.km <- clusterByKmeans(data.train,3)
# Visualize the clustering results using plotcluster() from the fpc library
plotcluster(data.train, fit.km)
```



As can be seen from the above plot, there exist lots of overlaps between clusters, for instance, between cluster 2 and cluster 3. It is not suitable to say that the data are well-separated.

We can compare the clustering results by k-means and the original classification of the 178 data indicated by the variable "Type":

```
# A comparison of the clustering results by k-means
# and the original classification indicated by "Type"
wt.km <- table(wine$Type, fit.km)
wt.km
```

```
##      fit.km
##      1  2  3
##  1  0  1 58
```

```
## 2 25 20 26
## 3 2 28 18
```

```
# Calculate the adjusted rand index, which is used for
# quantifying to what extent the two ways of clustering agree with each other.
randIndex(wt.km)
```

```
## ARI
## 0.1984624
```

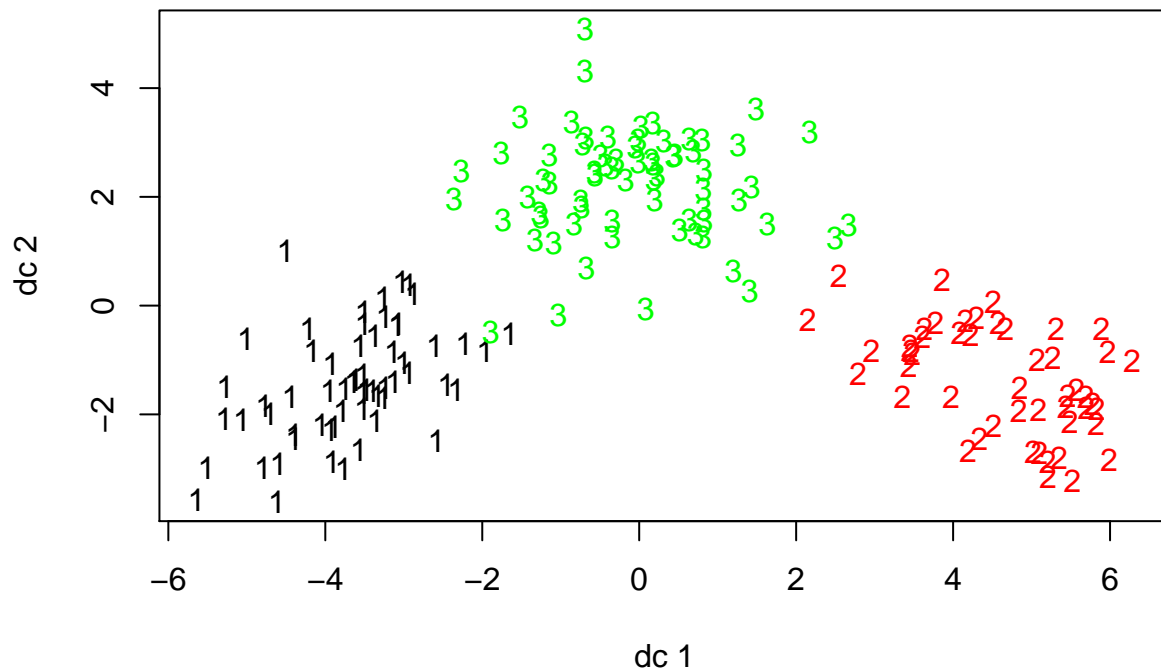
The above table shows directly that differences between the two indeed exist. The value of the adjusted rand index is 0.1984624, which is not quite close to 1. Note that the adjusted rand index which is a measure of the degree of agreement between two partitions takes values between -1 (no agreement) and 1 (perfect agreement). Hence, the partition obtained by K means does not match the original types of the wine very well.

Now we repeat the above exercise using scaled data:

```
# Obtain the wine data from the package "rattle"
data(wine, package="rattle")
head(wine)
```

```
## Type Alcohol Malic Ash Alkalinity Magnesium Phenols Flavanoids
## 1 1 14.23 1.71 2.43 15.6 127 2.80 3.06
## 2 1 13.20 1.78 2.14 11.2 100 2.65 2.76
## 3 1 13.16 2.36 2.67 18.6 101 2.80 3.24
## 4 1 14.37 1.95 2.50 16.8 113 3.85 3.49
## 5 1 13.24 2.59 2.87 21.0 118 2.80 2.69
## 6 1 14.20 1.76 2.45 15.2 112 3.27 3.39
## Nonflavanoids Proanthocyanins Color Hue Dilution Proline
## 1 0.28 2.29 5.64 1.04 3.92 1065
## 2 0.26 1.28 4.38 1.05 3.40 1050
## 3 0.30 2.81 5.68 1.03 3.17 1185
## 4 0.24 2.18 7.80 0.86 3.45 1480
## 5 0.39 1.82 4.32 1.04 2.93 735
## 6 0.34 1.97 6.75 1.05 2.85 1450
```

```
# Exclude the "Type" variable from the data inputs,
# center and scale the rest of the data and assign the rest to data.train
data.train <- scale(wine[,-1])
# Use set seed() so that the clustering results can be reproducible.
# set.seed(37810)
# Use k-means method to cluster the wines into 3 groups
fit.km <- clusterByKmeans(data.train,3)
# Visualize the clustering results using plotcluster() from the fpc library
plotcluster(data.train, fit.km)
```



The command, `data.train <- scale(wine[-1])`, is used for center and scale the wine data excluding the variable Type. More specifically, for each column of `wine[-1]`, first subtract the corresponding column mean and then dividing that centered column by its standard deviation. In `data.train`, each column has a mean of zero.

Using the scaled data, we obtain a better clustering result in the sense that now the clusters seem to be well-separated, despite a minor overlap between cluster 1 and cluster 3.

We can compare the clustering results by k-means and the original classification of the 178 data indicated by the variable “Type”:

```
# A comparison of the clustering results by k-means
# and the original classification indicated by "Type"
wt.km <- table(wine$Type, fit.km)
wt.km
```

```
##      fit.km
##      1  2  3
##  1 57  0  2
##  2  3  1 67
##  3  0 48  0
```

```
randIndex(wt.km)
```

```
##      ARI
## 0.8951549
```

Now the table shows that there are fewer missing samples than in the unscaled case. The value of the adjusted rand index is 0.6873931, which is much larger than in the unscaled case. Hence, scaling helps improve the performance of K means in the sense that after scaling the clustering result using K means is more similar to the original classification.

In the following we show repeat the above exercise for the iris dataset. We first use the original dataset which is not scaled.

```
# Obtain the iris data
```

```
data(iris)
```

```
head(iris)
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2   setosa
## 2         4.9         3.0         1.4         0.2   setosa
## 3         4.7         3.2         1.3         0.2   setosa
## 4         4.6         3.1         1.5         0.2   setosa
## 5         5.0         3.6         1.4         0.2   setosa
## 6         5.4         3.9         1.7         0.4   setosa
```

```
# Exclude the "Species" variable and assign the rest to data.train
```

```
data.train <- iris[-5]
```

```
# Use set seed() so that the clustering results can be reproducible
```

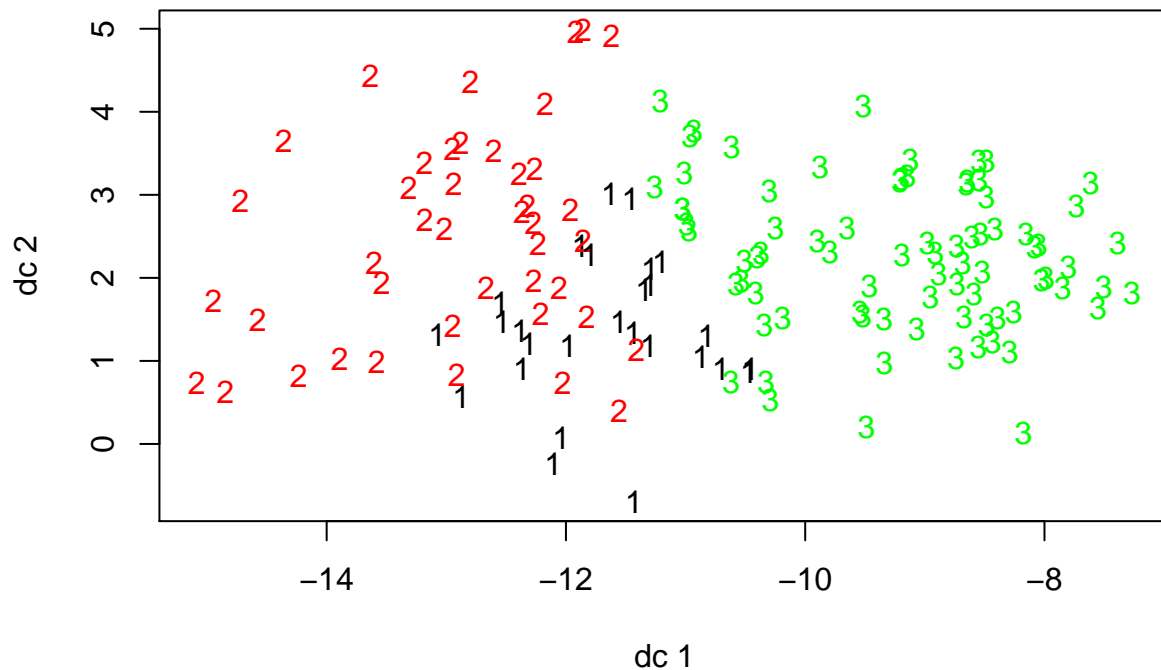
```
set.seed(37810)
```

```
# Use k-means method to cluster the data into 3 groups
```

```
fit.km <- clusterByKmeans(data.train,3)
```

```
# Visualize the clustering results using plotcluster() from the fpc library
```

```
plotcluster(data.train, fit.km)
```



The plot here shows some overlaps, in particular between cluster 1 and cluster 2 and between cluster 1 and cluster 3, as a result of which k-means does not work very well for clustering this dataset.

We can compare the clustering results by k-means and the original classification of the data indicated by the variable “species”:

```
fs.km <- table(iris$Species,fit.km)
fs.km
```

```
##          fit.km
##          1  2  3
## setosa      0  0 50
## versicolor 25  3 22
## virginica   2 40  8
```

```
randIndex(fs.km)
```

```
##          ARI
## 0.4305639
```

The adjusted rand index here is 0.4305639, which is not quite close to 1. So the algorithm’s clusters do not match the original partition very well.

Now we instead use the scaled iris dataset:


```
# Obtain the iris data
```

```
data(iris)
```

```
head(iris)
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2   setosa
## 2         4.9         3.0         1.4         0.2   setosa
## 3         4.7         3.2         1.3         0.2   setosa
## 4         4.6         3.1         1.5         0.2   setosa
## 5         5.0         3.6         1.4         0.2   setosa
## 6         5.4         3.9         1.7         0.4   setosa
```

```
# Exclude the "Species" variable and assign the rest to data.train
```

```
data.train <- scale(iris[1:5])
```

```
# Use set seed() so that the clustering results can be reproducible
```

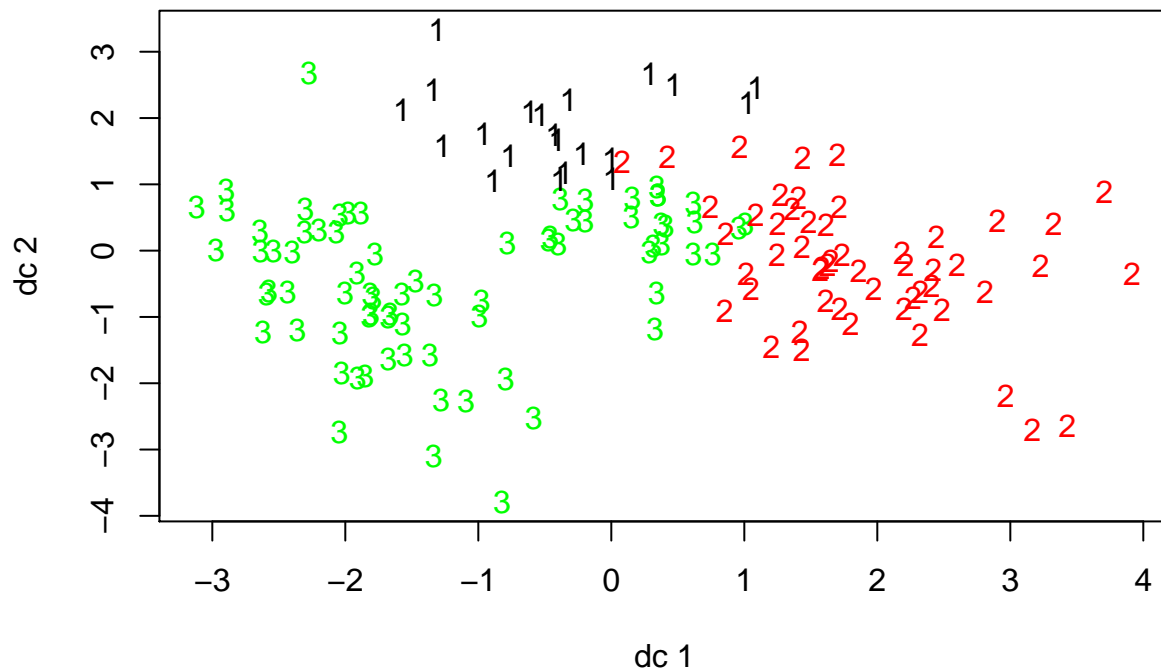
```
set.seed(37810)
```

```
# Use k-means method to cluster the data into 3 groups
```

```
fit.km <- clusterByKmeans(data.train,3)
```

```
# Visualize the clustering results using plotcluster() from the fpc library
```

```
plotcluster(data.train, fit.km)
```



As shown in the plot, there are still some overlaps between clusters, indicating that k-means does not work very well even for this scaled data set.

To quantify the comparison, we can still apply `randIndex()`

```
fs.km <- table(iris$Species,fit.km)
fs.km
```

```
##           fit.km
##           1  2  3
## setosa      0  0 50
## versicolor 19 12 19
## virginica   2 41  7
```

```
randIndex(fs.km)
```

```
##           ARI
## 0.3910336
```

The adjusted rand index now is 0.3910336, which is even a little bit smaller than the unscaled case. In other words, for the iris dataset, scaling does not help improving the effectiveness of the k-means.

As a summary, scaling is useful for improving the performance of K means for the wine data but not for the iris data.